```python
In [1]:  import pandas as pd
         import numpy as np
         import matplotlib.pyplot as plt
         from sklearn.preprocessing import MinMaxScaler,OneHotEncoder,LabelEncoder,StandardS
         from sklearn.model_selection import train_test_split
         from sklearn.neighbors import KNeighborsClassifier
         from sklearn.model_selection import GridSearchCV
         from sklearn.metrics import accuracy_score,precision_score,recall_score,f1_score,co
         from sklearn.model_selection import KFold,cross_validate
```

# Read data

```python
In [2]:  col_names = ['fLength', 'fWidth', 'fSize', 'fConc', 'fConc1', 'fAsym', 'fM3Long', '
         df = pd.read_csv("magic04.data")
         df.columns = col_names
         df.head()
```

Out[2]:

| | fLength | fWidth | fSize | fConc | fConc1 | fAsym | fM3Long | fM3Trans | fAlpha | fl |
|---|---------|--------|-------|-------|--------|-------|---------|----------|--------|----|
| 0 | 31.6036 | 11.7235 | 2.5185 | 0.5303 | 0.3773 | 26.2722 | 23.8238 | -9.9574 | 6.3609 | 205. |
| 1 | 162.0520 | 136.0310 | 4.0612 | 0.0374 | 0.0187 | 116.7410 | -64.8580 | -45.2160 | 76.9600 | 256. |
| 2 | 23.8172 | 9.5728 | 2.3385 | 0.6147 | 0.3922 | 27.2107 | -6.4633 | -7.1513 | 10.4490 | 116. |
| 3 | 75.1362 | 30.9205 | 3.1611 | 0.3168 | 0.1832 | -5.5277 | 28.5525 | 21.8393 | 4.6480 | 356. |
| 4 | 51.6240 | 21.1502 | 2.9085 | 0.2420 | 0.1340 | 50.8761 | 43.1887 | 9.8145 | 3.6130 | 238. |

```python
In [3]:  print(set(df['class']))
```

```
{'g', 'h'}
```

```python
In [4]:  df.dropna(inplace=True)
```

```python
In [5]:  plt.bar(['gamma','hadron'],height=[len(df[df['class']=='g']),len(df[df['class']=='h
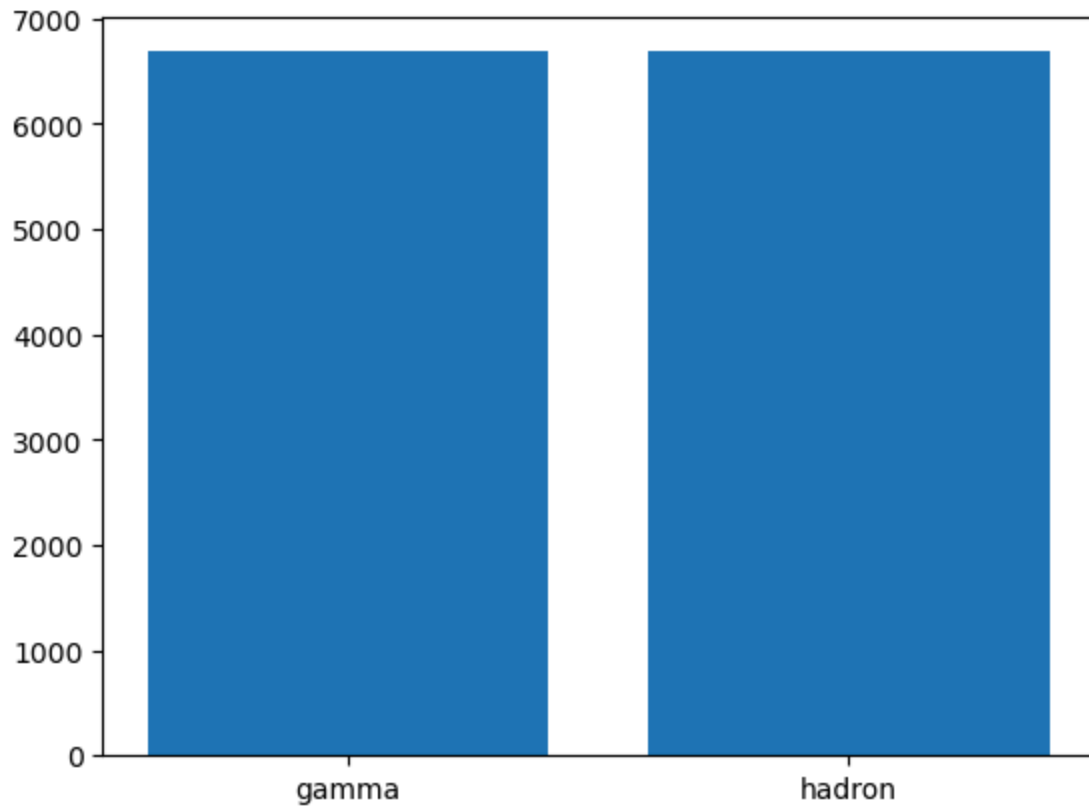```

```
Out[5]:  <BarContainer object of 2 artists>
```

## Balancing Data

```
In [6]:  df = df.groupby('class').sample(len(df[df['class']=='h']))
```
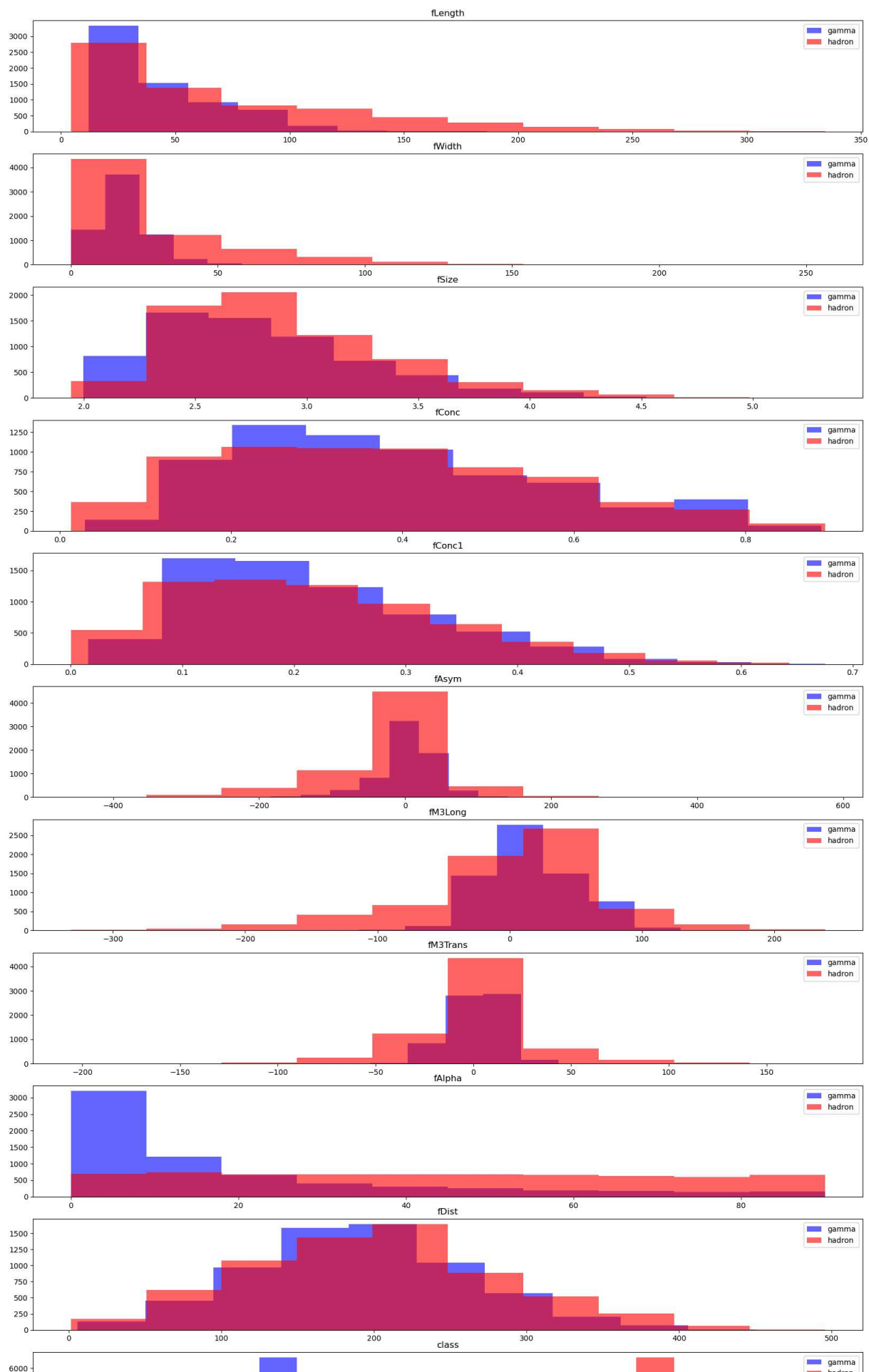
```
In [7]:  plt.bar(['gamma','hadron'],height=[len(df[df['class']=='g']),len(df[df['class']=='h
```

Out[7]:  `<BarContainer object of 2 artists>`

# Visualizing Data Columns

```
In [8]:  fig,axs = plt.subplots(11,figsize=(20,35))
         for i,label in enumerate(df):
             axs[i].hist(df[df['class']=='g'][label],color='blue',alpha=0.6,label="gamma")
             axs[i].hist(df[df['class']=='h'][label],color='red',alpha=0.6,label="hadron")
             axs[i].title.set_text(label)
             axs[i].legend()
         plt.show()
```

# Encoding

```
In [9]:  en = LabelEncoder()
         df['class'] = en.fit_transform(df['class'])
```

```
In [10]:  df.head()
```

Out[10]:

| | fLength | fWidth | fSize | fConc | fConc1 | fAsym | fM3Long | fM3Trans | fAlpha | |
|---|---|---|---|---|---|---|---|---|---|---|
| 4710 | 41.7854 | 17.7204 | 2.6365 | 0.3210 | 0.1721 | -17.1632 | -37.2658 | 13.6465 | 0.5310 | 1 |
| 10250 | 41.8464 | 19.4302 | 3.1166 | 0.2217 | 0.1235 | 11.5249 | 27.9829 | -11.9406 | 0.0920 | 1 |
| 3165 | 95.2562 | 27.9901 | 3.1992 | 0.2851 | 0.1489 | -66.8185 | -71.2218 | -24.3355 | 3.6810 | 3 |
| 2855 | 13.0339 | 11.1611 | 2.0810 | 0.7552 | 0.3942 | 14.9112 | -4.3689 | 11.7397 | 86.3258 | |
| 4526 | 48.2713 | 11.7368 | 2.5192 | 0.5446 | 0.2950 | 12.0288 | 40.4022 | 1.4286 | 10.3490 | 2 |

```
In [11]:  import seaborn as sns
          correlation_matrix = df.corr()

          plt.figure(figsize=(10, 8))

          sns.heatmap(correlation_matrix, annot=True, cmap="coolwarm", vmin=-1, vmax=1, linew

          plt.title("Correlation Heatmap")

          plt.show()
```

Correlation Heatmap

# splitting

```
In [12]:  X = df.iloc[:,:-1]
          y = df.iloc[:,-1]
```

```
In [13]:  norm = StandardScaler()
          X = norm.fit_transform(X)
```

```
In [14]:  X_train, X_test, y_train, y_test = train_test_split(X,y,random_state=42,shuffle=Tru
          X_val, X_test, y_val, y_test = train_test_split(X_test,y_test,random_state=42,shuff
```

```
In [15]:  plt.bar(['train','validation','test'],[len(y_train),len(y_val),len(y_test)])
```

```
Out[15]:  <BarContainer object of 3 artists>
```

In [16]: 
```python
neighbors = range(1,100,2)
```

In [17]: 
```python
def compare(neighbors, X_train, X_val, y_train, y_val):
    accs = []
    prec = []
    rec = []
    f1_scores = []
    CMs = []

    for i in neighbors:
        model = KNeighborsClassifier(n_neighbors=i)

        # Train the model on the full training data
        model.fit(X_train, y_train)

        # Predict on training and validation sets
        y_train_pred = model.predict(X_train)
        y_val_pred = model.predict(X_val)

        # Calculate metrics for training set
        train_acc = accuracy_score(y_train, y_train_pred)
        train_precision = precision_score(y_train, y_train_pred)
        train_recall = recall_score(y_train, y_train_pred)
        train_f1 = f1_score(y_train, y_train_pred)

        # Calculate metrics for validation set
        val_acc = accuracy_score(y_val, y_val_pred)
        val_precision = precision_score(y_val, y_val_pred)
        val_recall = recall_score(y_val, y_val_pred)
        val_f1 = f1_score(y_val, y_val_pred)
```

```
            # Confusion matrix for validation
            cm = confusion_matrix(y_val, y_val_pred)

            # Append training and validation metrics
            accs.append((val_acc,train_acc))
            prec.append((val_precision,train_precision))
            rec.append((val_recall,train_recall))
            f1_scores.append((val_f1,train_f1))
            CMs.append(cm)

    return accs, prec, rec, f1_scores, CMs
```

In [27]:
```
accs,prec,rec,f1_scores,CMs = compare(neighbors,X_train,X_val,y_train,y_val)
```

In [28]:
```python
val_accs = [x[0] for x in accs]
train_accs = [x[1] for x in accs]
val_prec = [x[0] for x in prec]
train_prec = [x[1] for x in prec]
val_rec = [x[0] for x in rec]
train_rec = [x[1] for x in rec]
val_f1 = [x[0] for x in f1_scores]
train_f1 = [x[1] for x in f1_scores]

# Plot the results
plt.figure(figsize=(14, 10))

# Plot accuracy
plt.plot(neighbors, val_accs, label='Validation Accuracy', marker='o', color='blue'

# Plot precision
plt.plot(neighbors, val_prec, label='Validation Precision', marker='o', color='gree

# Plot recall
plt.plot(neighbors, val_rec, label='Validation Recall', marker='o', color='orange')

# Plot F1 score
plt.plot(neighbors, val_f1, label='Validation F1 Score', marker='o', color='red')

# Add labels, title, and legend
plt.title('K-Neighbors Classifier Performance with Varying n_neighbors')
plt.xlabel('Number of Neighbors (n_neighbors)')
plt.ylabel('Score')
plt.legend()
plt.grid(True)
plt.show()
```
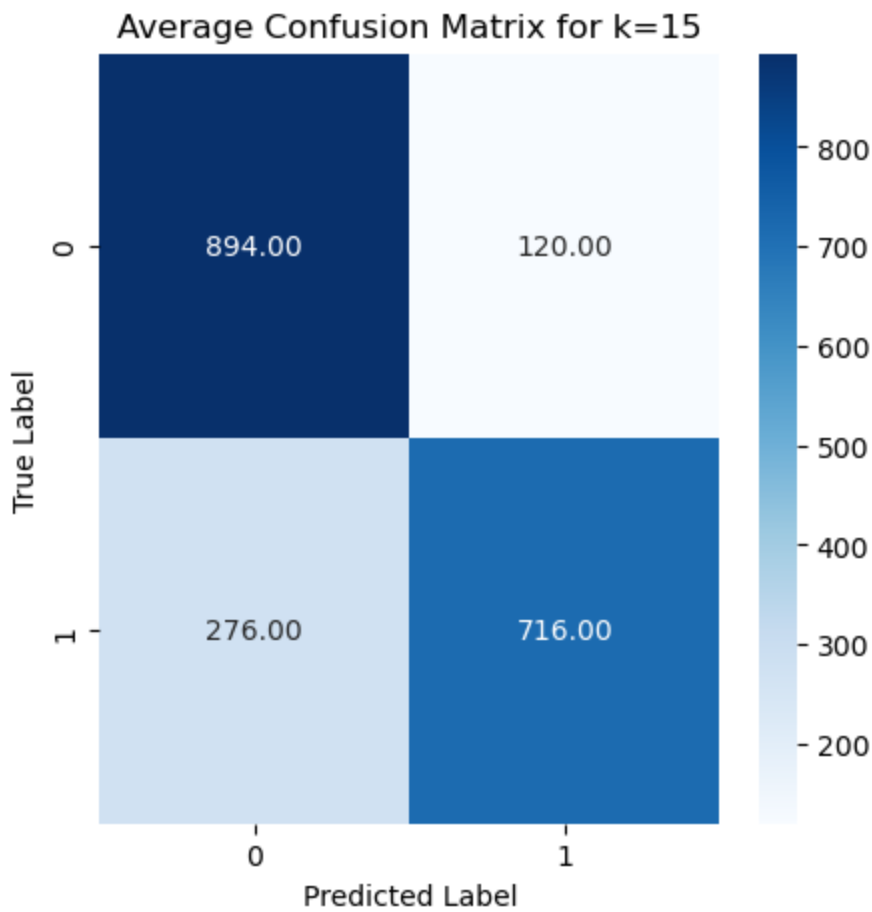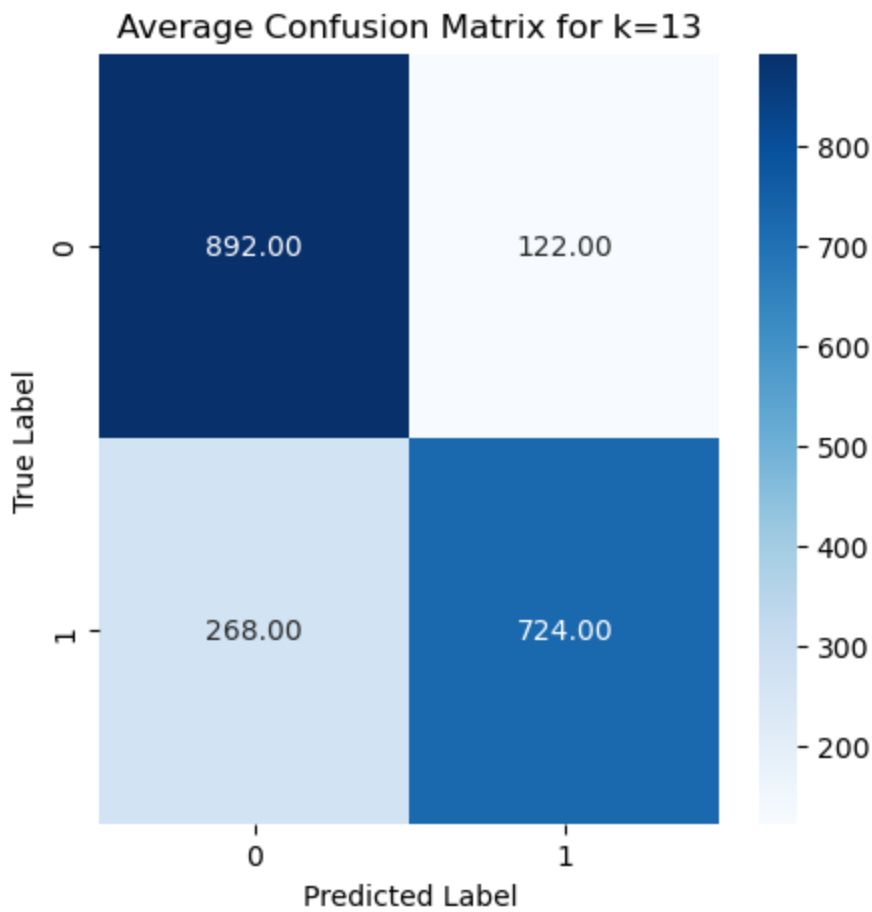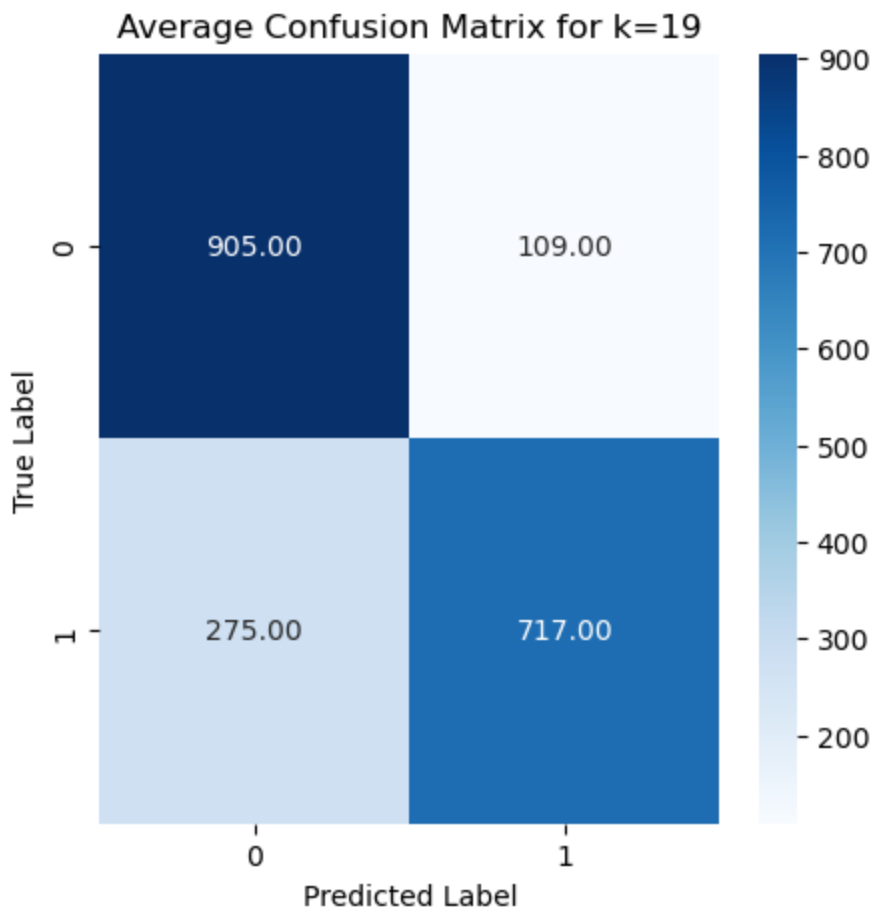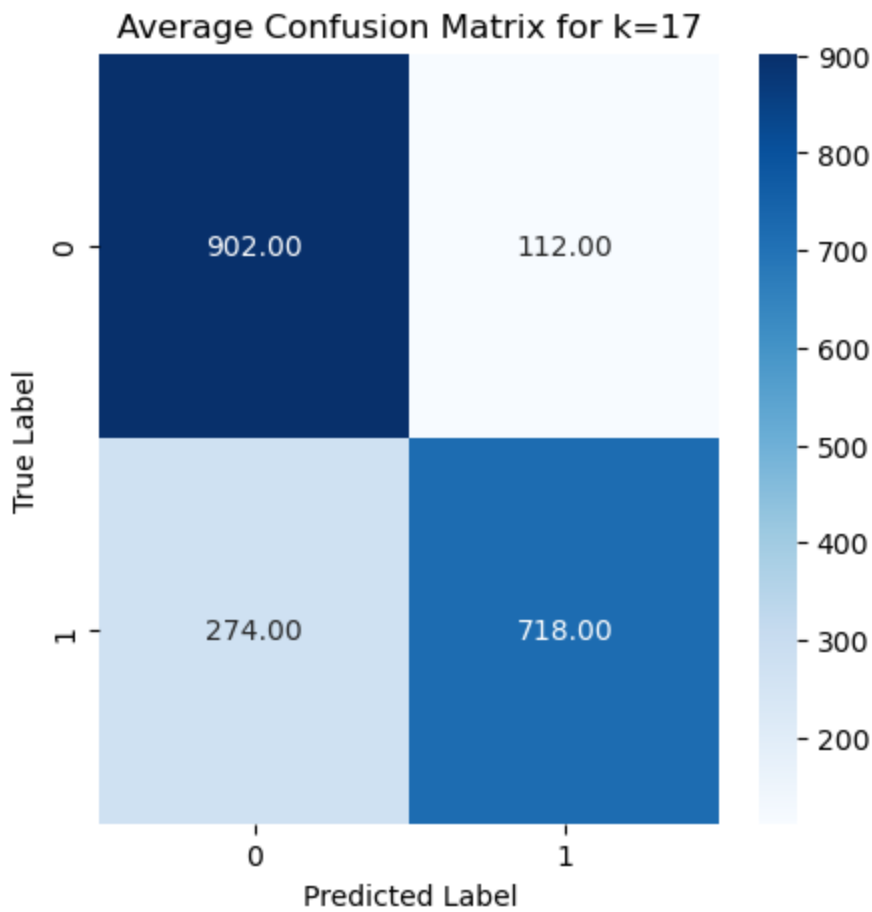
K-Neighbors Classifier Performance with Varying n_neighbors

```
In [29]:  plt.figure(figsize=(14, 10))
          plt.plot(neighbors, train_accs, label='Training Accuracy', marker='o', linestyle='-
          plt.plot(neighbors, train_prec, label='Training Precision', marker='o', linestyle='
          plt.plot(neighbors, train_rec, label='Training Recall', marker='o', linestyle='--',
          plt.plot(neighbors, train_f1, label='Training F1 Score', marker='o', linestyle='--'
          plt.title('K-Neighbors Classifier Performance with Varying n_neighbors')
          plt.xlabel('Number of Neighbors (n_neighbors)')
          plt.ylabel('Score')
          plt.legend()
          plt.grid(True)
          plt.show()
```

K-Neighbors Classifier Performance with Varying n_neighbors



In [30]:
```python
import seaborn as sns
for i, cm in enumerate(CMs):
    plt.figure(figsize=(5, 5))
    sns.heatmap(cm, annot=True, fmt='.2f', cmap='Blues')
    plt.title(f'Average Confusion Matrix for k={neighbors[i]}')
    plt.ylabel('True Label')
    plt.xlabel('Predicted Label')
    plt.show()
```

## Average Confusion Matrix for k=1



## Average Confusion Matrix for k=3

Average Confusion Matrix for k=5



Average Confusion Matrix for k=7

## Average Confusion Matrix for k=9



## Average Confusion Matrix for k=11

## Average Confusion Matrix for k=13



## Average Confusion Matrix for k=15

## Average Confusion Matrix for k=17



## Average Confusion Matrix for k=19

## Average Confusion Matrix for k=21



## Average Confusion Matrix for k=23

Average Confusion Matrix for k=25



Average Confusion Matrix for k=27

## Average Confusion Matrix for k=29



## Average Confusion Matrix for k=31

## Average Confusion Matrix for k=33



## Average Confusion Matrix for k=35

## Average Confusion Matrix for k=37



## Average Confusion Matrix for k=39

## Average Confusion Matrix for k=41



## Average Confusion Matrix for k=43

## Average Confusion Matrix for k=45



## Average Confusion Matrix for k=47

## Average Confusion Matrix for k=49



## Average Confusion Matrix for k=51

## Average Confusion Matrix for k=53



## Average Confusion Matrix for k=55

## Average Confusion Matrix for k=57



## Average Confusion Matrix for k=59

## Average Confusion Matrix for k=61



## Average Confusion Matrix for k=63

## Average Confusion Matrix for k=65



## Average Confusion Matrix for k=67

## Average Confusion Matrix for k=69



## Average Confusion Matrix for k=71

## Average Confusion Matrix for k=73



## Average Confusion Matrix for k=75

## Average Confusion Matrix for k=77



## Average Confusion Matrix for k=79

## Average Confusion Matrix for k=81



## Average Confusion Matrix for k=83

## Average Confusion Matrix for k=85



## Average Confusion Matrix for k=87

## Average Confusion Matrix for k=89



## Average Confusion Matrix for k=91

## Average Confusion Matrix for k=93



## Average Confusion Matrix for k=95

Average Confusion Matrix for k=97



Average Confusion Matrix for k=99

```
In [31]:  best_k = neighbors[np.argmax(val_accs)]

          print(f"best K value = {best_k}")
          print(f"accuracy = {val_accs[np.argmax(val_accs)]}")
          print(f"precision = {val_prec[np.argmax(val_accs)]}")
          print(f"Recall = {val_rec[np.argmax(val_accs)]}")
          print(f"F1 = {val_f1[np.argmax(val_accs)]}")
```
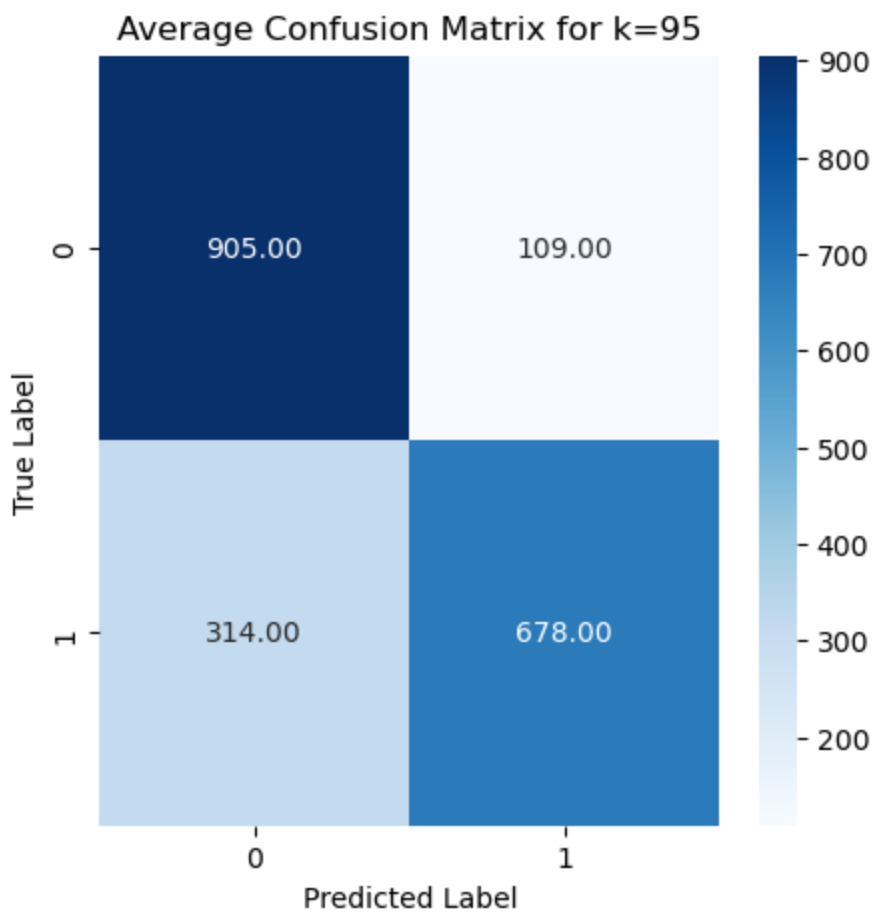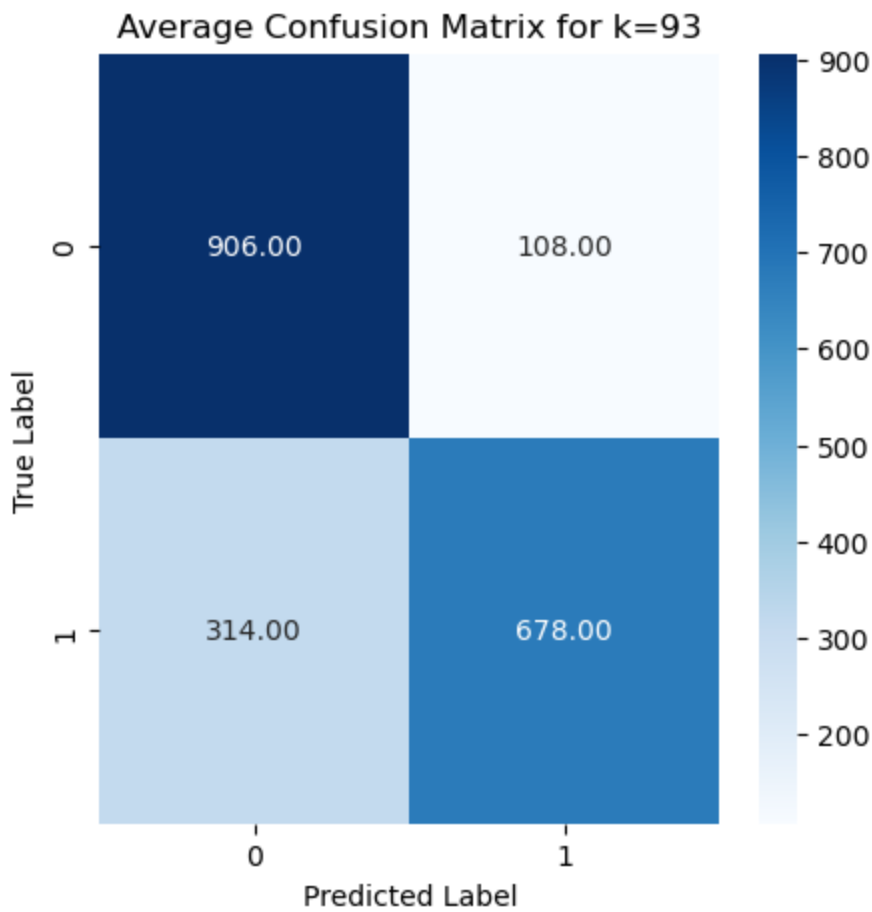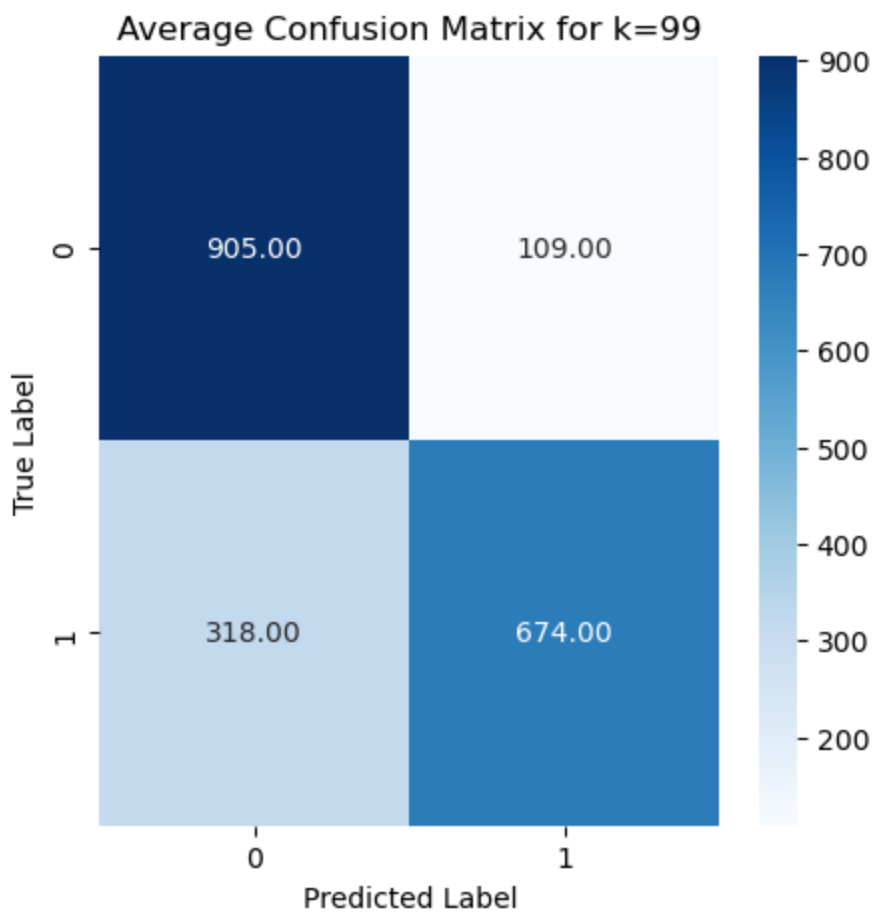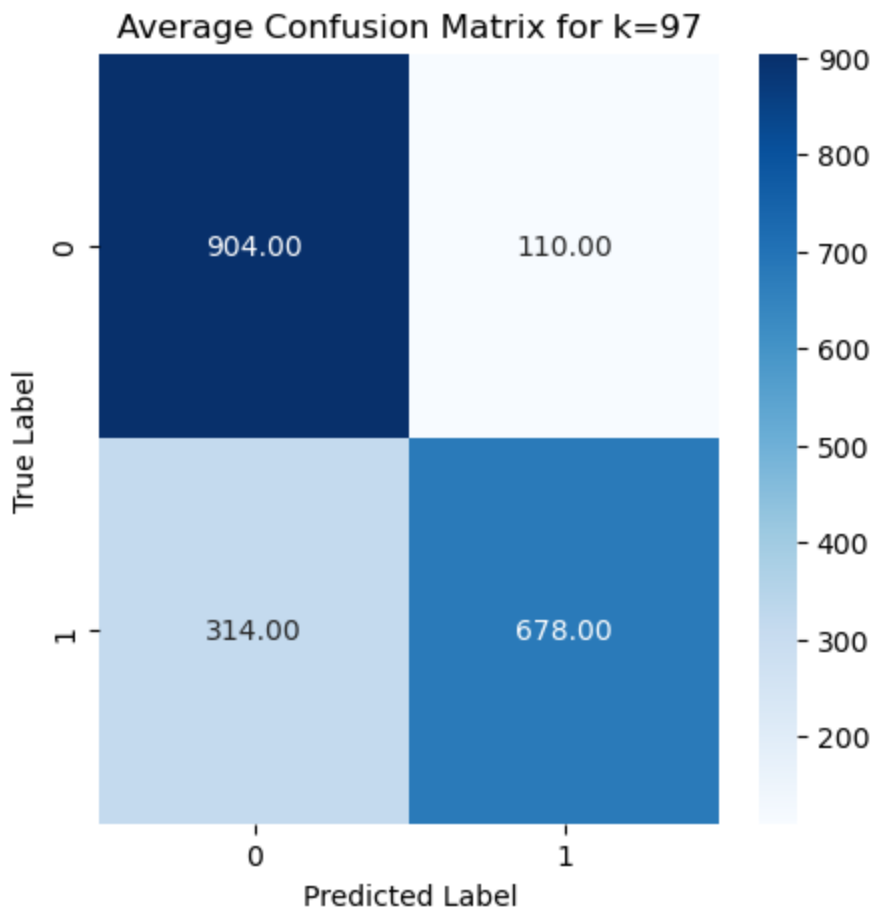
```
best K value = 9
accuracy = 0.8155533399800599
precision = 0.8624708624708625
Recall = 0.7459677419354839
F1 = 0.8
```

# Testing

```
In [32]:  model = KNeighborsClassifier(n_neighbors=best_k)
          model.fit(X_train,y_train)
```
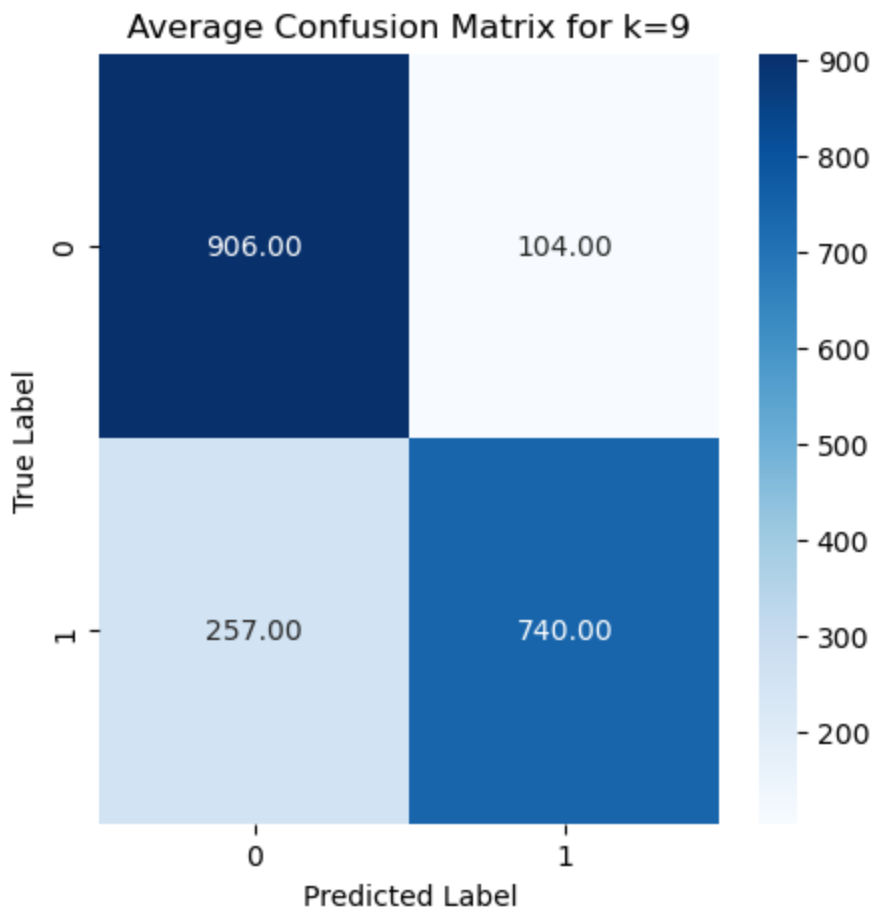
```
Out[32]:  ▼      KNeighborsClassifier      ⓘ ⑦

          KNeighborsClassifier(n_neighbors=9)
```

```
In [33]:  y_pred = model.predict(X_test)
```

```
In [34]:  acc = accuracy_score(y_test,y_pred)
          prec = precision_score(y_test,y_pred)
          rec = recall_score(y_test,y_pred)
          f1 = f1_score(y_test,y_pred)
          print(f"accuracy = {round(acc,2)}")
          print(f"precision = {round(prec,2)}")
          print(f"Recall = {round(rec,2)}")
          print(f"F1 = {round(f1,2)}")
```

```
accuracy = 0.82
precision = 0.88
Recall = 0.74
F1 = 0.8
```

```
In [35]:  cm = confusion_matrix(y_test, y_pred)
          plt.figure(figsize=(5, 5))
          sns.heatmap(cm, annot=True, fmt='.2f', cmap='Blues')
          plt.title(f'Average Confusion Matrix for k={best_k}')
          plt.ylabel('True Label')
          plt.xlabel('Predicted Label')
          plt.show()
```

Average Confusion Matrix for k=9

In [ ]: